

# Mathematics Performance in Secondary Education

Yao-Chih Hsu, Xuan-Chun Wang, and Wen-Lee Sin

June 5, 2025

# Outline

- 1 Introduction
- 2 Data Visualization
- 3 Dimensionality Reduction
- 4 Regression
- 5 Conclusion
- 6 References

# Introduction

- Dataset: *Math-Students Performance Data* from Kaggle (Shamim, 2025).

# Introduction

- Dataset: *Math-Students Performance Data* from Kaggle (Shamim, 2025).
- Variables G1, G2, G3, and absences were provided by the school.
- Remaining variables were collected via questionnaires and are mostly categorical.

# Goal

Analyze G3 to identify influential variables.

# Data Grouping

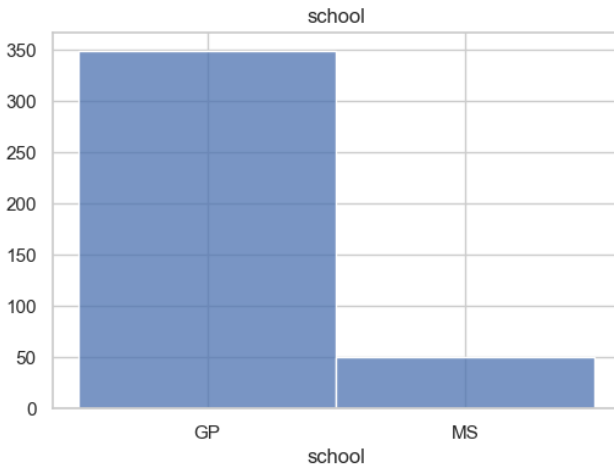
We grouped the data into the following:

Groups	Variables
support	schoolsup, famsup, paid
family	address, famsize, Pstatus, guardian, traveltime, famrel
parents	Medu, Fedu, Mjob, Fjob
performance	failures, studytime, absences
alcohol	Dalc, Walc, health
after_class	activities, freetime, goout
school_choice	reason, nursery, higher
score	G1, G2, G3

The variables not yet assigned to any group are:

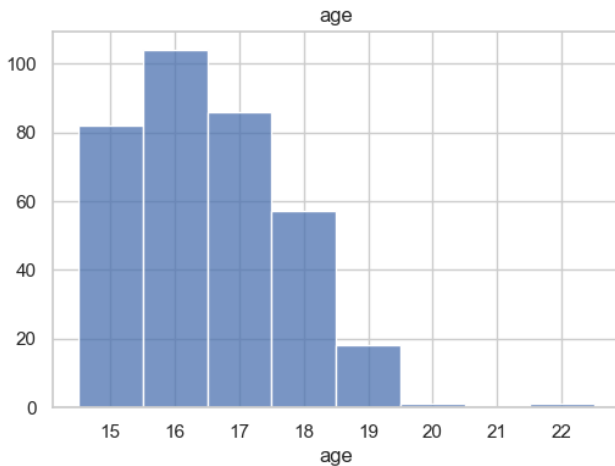
sex, age, internet, romantic.

# Data Visualization

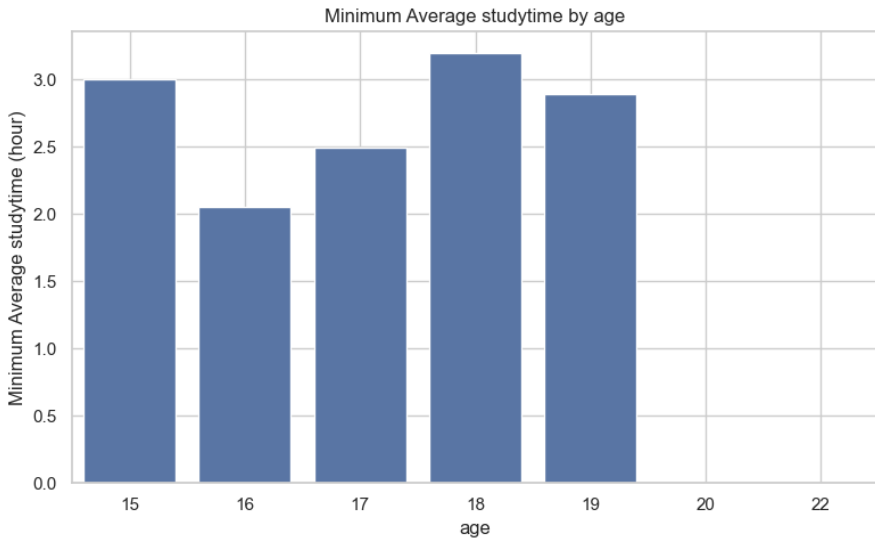


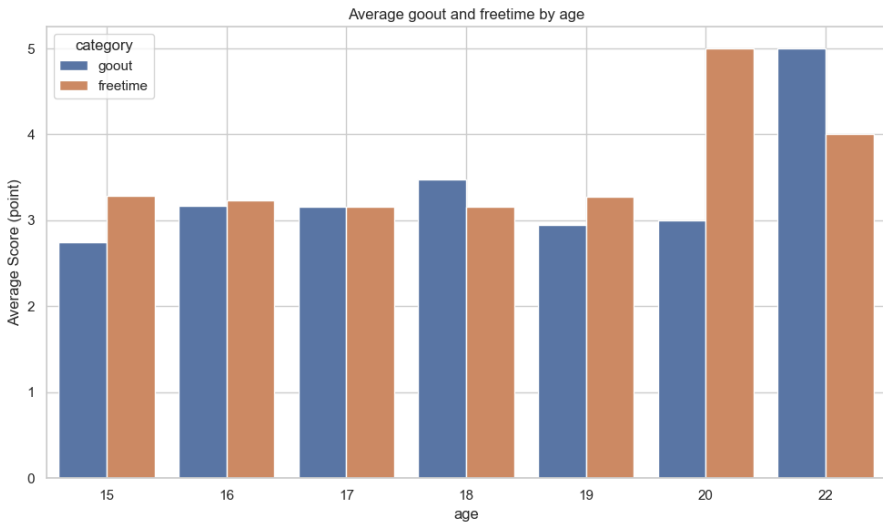
**GP (Gabriel Pereira)**

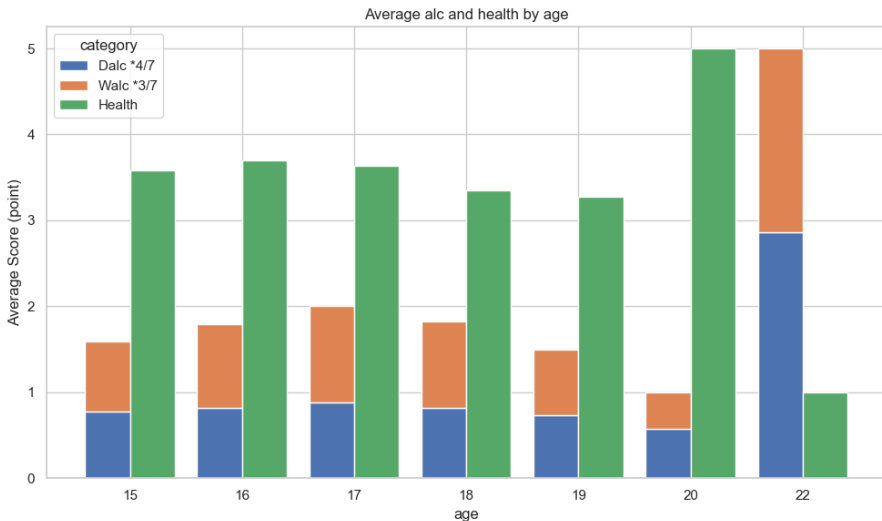
**MS (Mousinho da Silveira)**

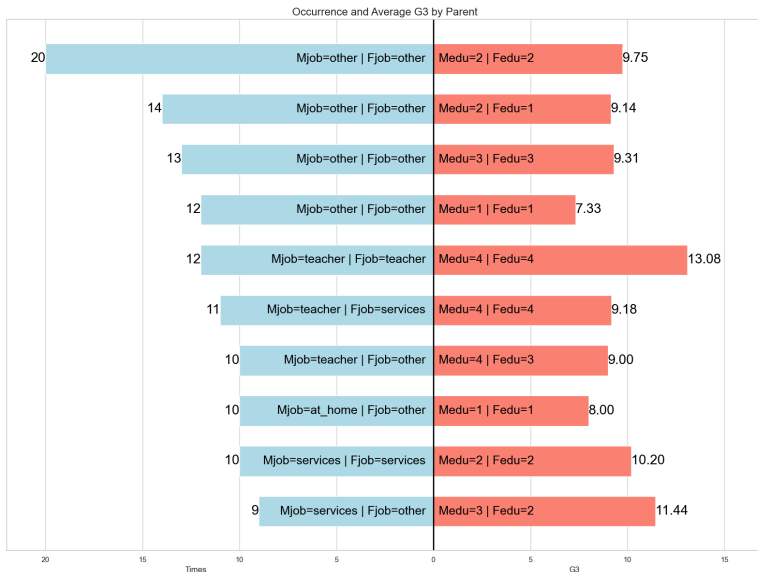


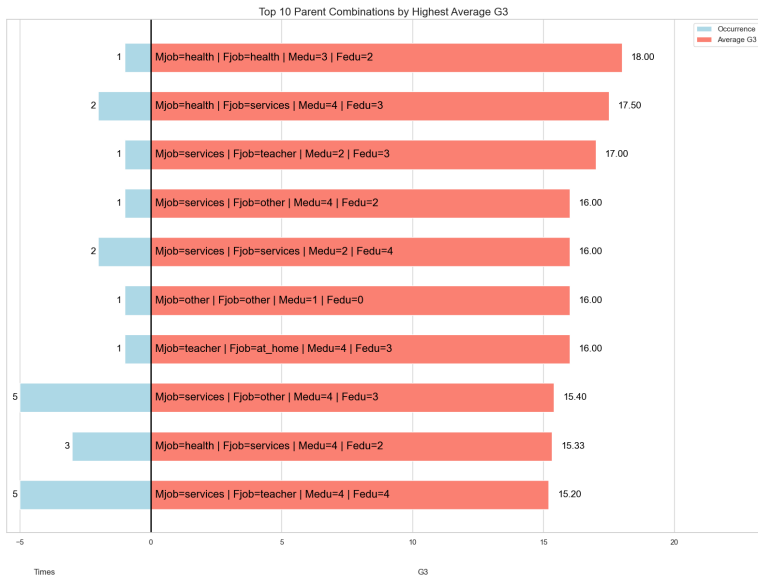




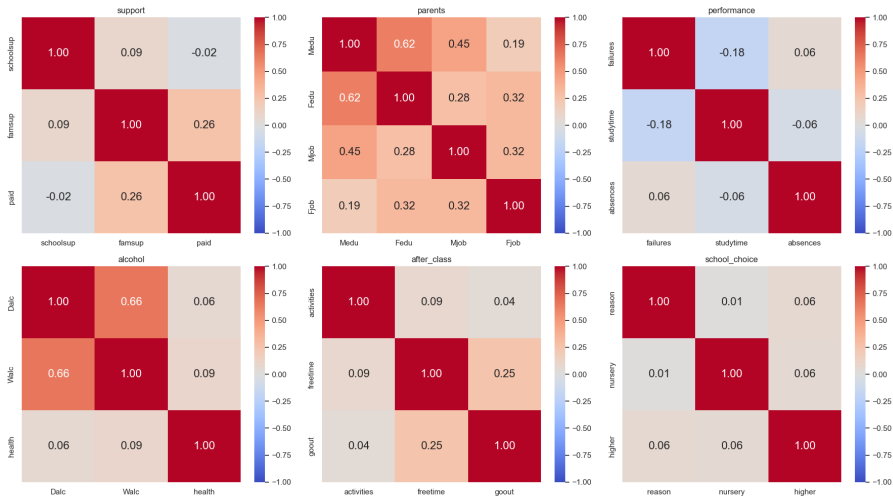




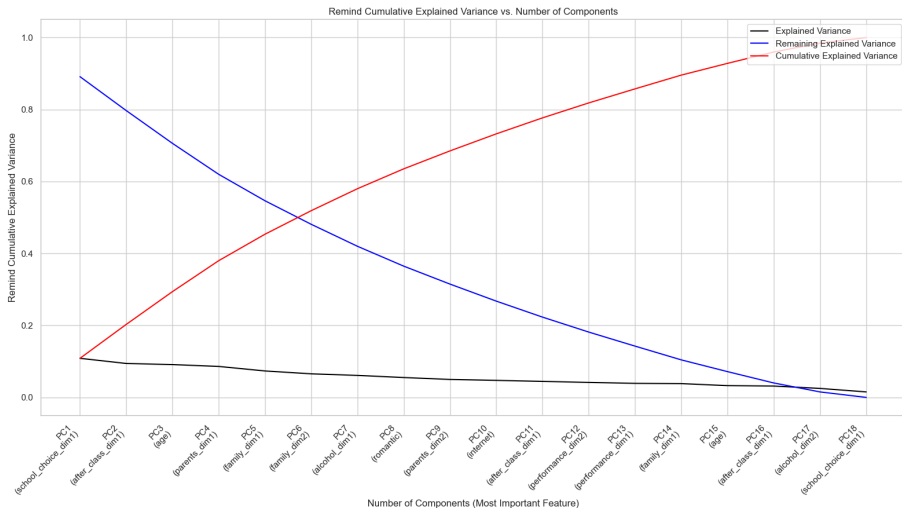




# Multidimensional Scaling



# PCA



# Variables of MLR

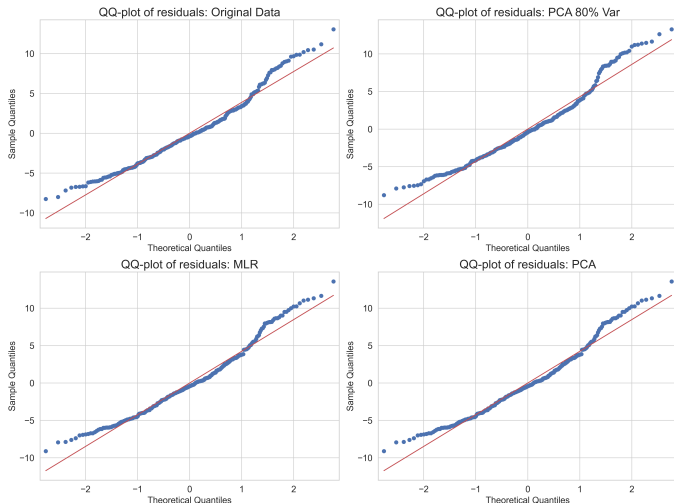
We compared four models:

- Original + MLR
- MDS + MLR
- MDS + PCA + MLR
- MDS + PCA 80% + MLR

The variables in the dataset transformed by MDS are as follows:



# Q-Q plots of MLR

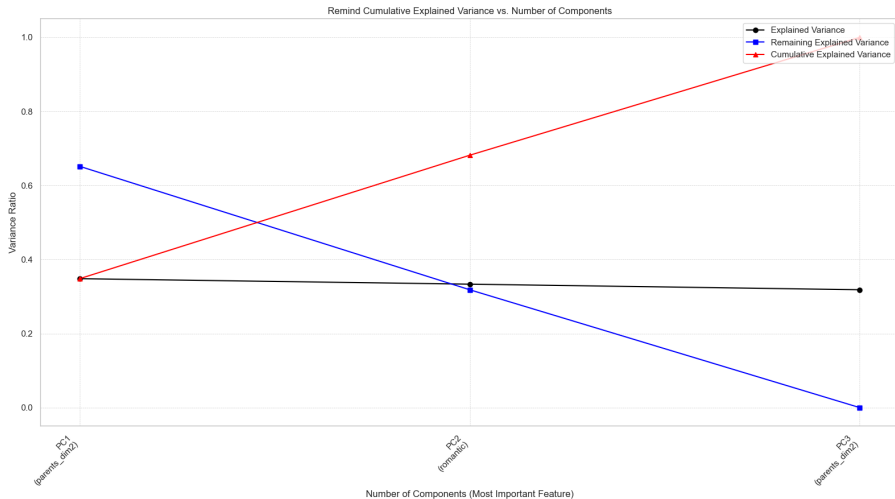


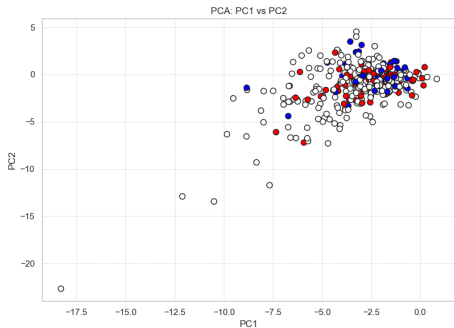
# Variables of Reduced MLR

Method	Selected Features
Original (Reduced) + MLR	failures, goout, Mjob_at_home, romantic, schoolsup
MDS (Reduced) + PCA + MLR	after_class_dim2, romantic, parents_dim2

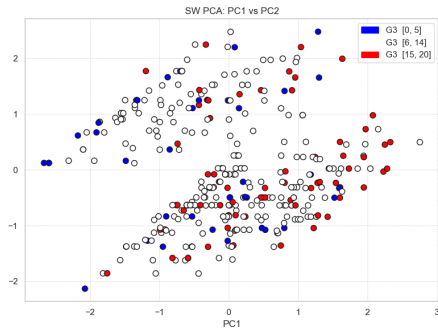
Table 1: Selected Features by Different Methods

# Original data + MDS + Stepwise + PCA





PC1 (school\_choice\_dim1)  
PC2 (after\_class\_dim1)



PC1 (parent\_dim2)  
PC2 (romantic)

# Conclusion

Model	$R^2$	Adj. $R^2$	MSE	P-value	A.I.C
Original + MLR	<b>0.30</b>	<b>0.20</b>	<b>58.03</b>	0.9985	2013.87
MDS + MLR	0.16	0.11	65.08	0.5948	2036.70
MDS + PCA + MLR	0.16	0.11	65.08	0.9639	2036.70
MDS + PCA 80% + MLR	0.13	0.10	81.21	0.9227	2035.49
Original (Reduced) + MLR	0.19	0.18	288.71	<b>0.0375</b>	<b>1995.23</b>
MDS (Reduced) + PCA + MLR	0.08	0.07	204.60	0.7878	2036.42
Paper Proposed + MLR	0.19	0.17	176.54	0.6199	2003.04

**Table 2:** Comparison of model performance metrics

The variables selected by the paper:  
absences , schoolsup , higher , failures , Mjob

# Conclusion

Model	TOP1	TOP2	TOP3
<b>Original + MLR</b>	<b>higher</b>	<b>Fjob_teacher</b>	<b>failures</b>
MDS + MLR	romantic	support_dim2	after_class_dim2
MDS + PCA + MLR	alcohol_dim1	after_class_dim1	age
MDS + PCA 80% + MLR	performance_dim2	alcohol_dim1	age
<b>Original (Reduced) + MLR</b>	<b>failures</b>	<b>Mjob_at_home</b>	<b>schoolsup</b>
MDS (Reduced) + PCA + MLR	parents_dim2	romantic	parents_dim2
<b>Paper Proposed + MLR</b>	<b>higher</b>	<b>failures</b>	<b>Mjob_health</b>

Table 3: Most important feature to each model

# Conclusion

The variables we selected are:

failures, Mjob, schoolsup

# References



Cortez, P. (2008). *Student Performance* [Dataset]. UCI Machine Learning Repository. <https://doi.org/10.24432/C5TG7T>



Cortez, P., & Silva, A. M. (2008). Using data mining to predict secondary school student performance.



Shamim, A. (2025). *Math students performance data*. Kaggle. <https://www.kaggle.com/datasets/adilshamim8/math-students>



Thanks for listening!